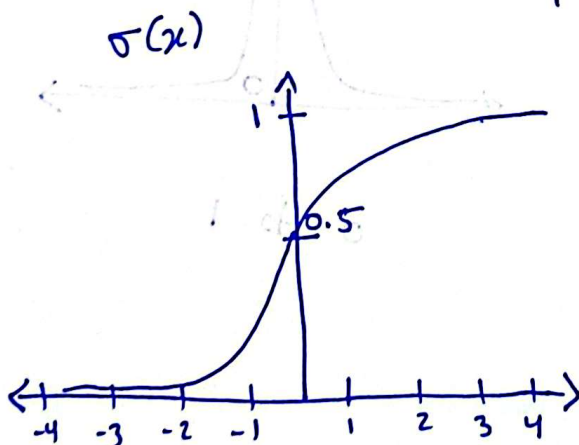# Activation Functions

Activation functions helps to determine the output of a neural network. These type of functions are attached to each neuron in the network, and determines whether it should be activated or not, based on whether each neuron's input is relevant for the model's prediction
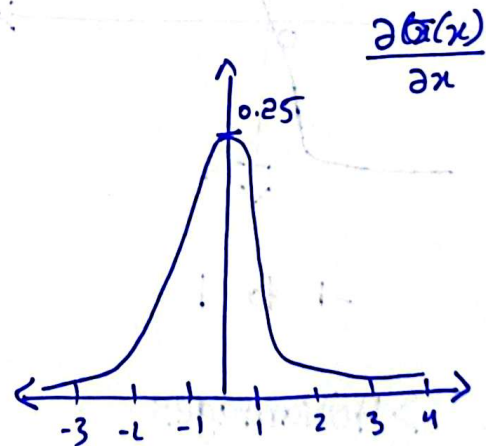
Activation functions also helps to normalize the output of each neuron to a specific range

$\longrightarrow$ Sigmoid Activation Function

$$\sigma(x) = \frac{1}{1+e^{-x}} \quad , \text{ range : 0 to } \underline{1}$$

$\sigma(x)$

$\frac{\partial \sigma(x)}{\partial x}$

0 to 1

0 to 0.25

> Advantages
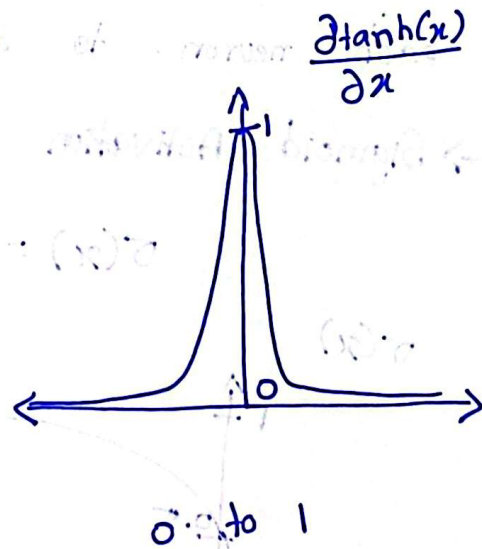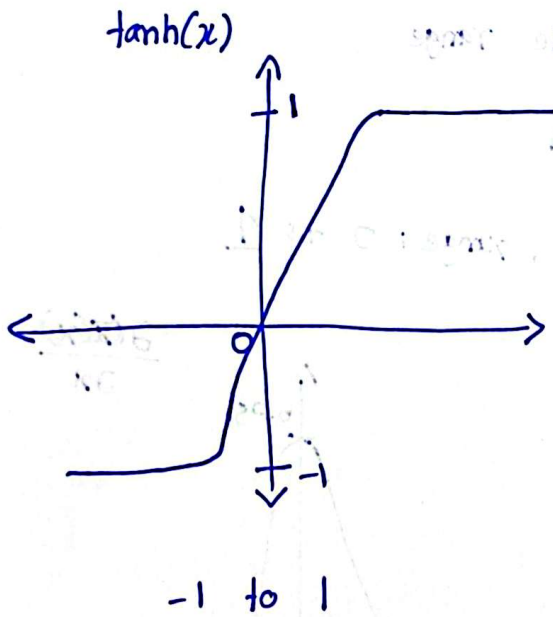
① Smooth gradient, preventing "jumps" in output values

② Output values bound between 0 and 1, normalizing the output of each neuron

③ Clear predictions, i.e. very close to 1 or 0.

> Disadvantages

① Prone to gradient vanishing

② Function output is not zero-centered

③ Power operations are relatively time consuming

——→ Tanh Activation Function

$$tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

tanh(x)

$\frac{\partial tanh(x)}{\partial x}$

−1 to 1

0 to 1

> Advantages

① Zero centric functions

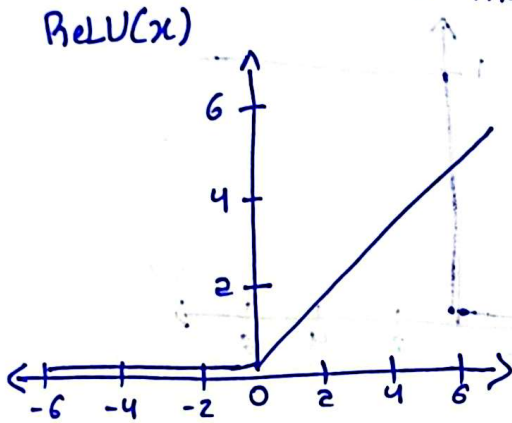② Modified version of sigmoid and comparatively better sigmoid function

> Disadvantages

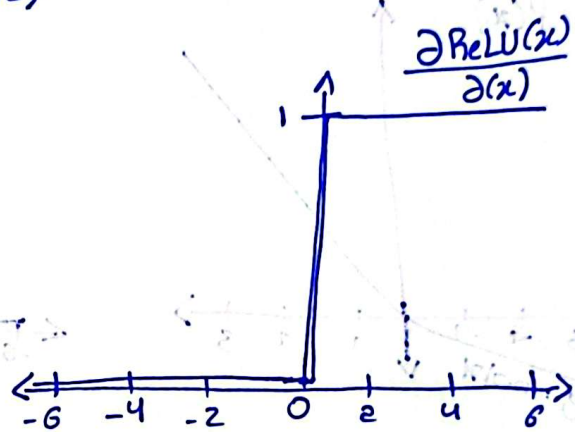① Vanishing gradient issue for large layers

② Expensive computational power

# → ReLU Activation Function

$$x = \begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } x \geq 0 \end{cases}$$

$$= max(0, x)$$

ReLU(x)

0 to ∞

$\dfrac{\partial ReLU(x)}{\partial(x)}$

0 or 1

if ReLU output is 1 => weight updation

if ReLU output is 0 => Dead neuron

> Advantages

① Solving vanishing gradient problem

② It doesnot activate all the neurons

③ Computation is very fast.

> Disadvantages
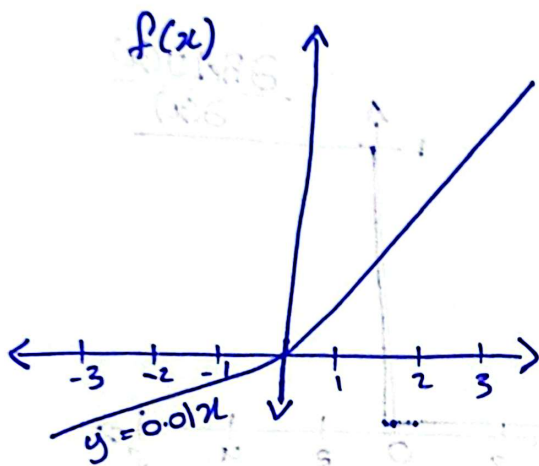
① ReLU function is not zero centric

② Dead neurons as it deactivates the neurons
   for -ve value and do not provide -ve direction.

③ For -ve input it does not do anything, which is not
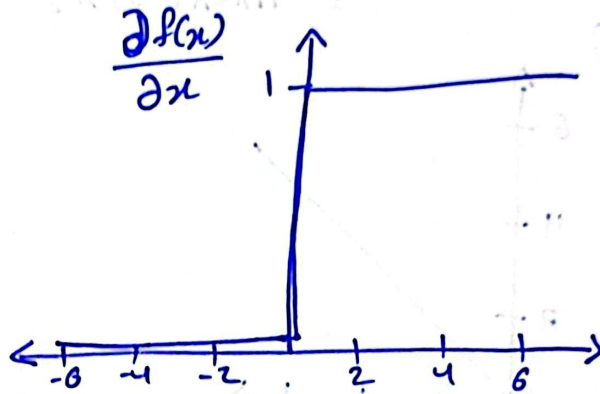   good for learning of Neural Network.

# —> Leaky ReLU

$$\begin{cases} 0.01x & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$$



$f(x)$

$y = 0.01x$

-3  -2  -1  1  2  3

$\dfrac{\partial f(x)}{\partial x}$
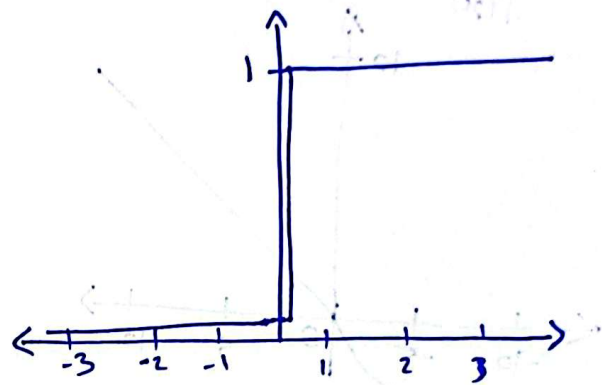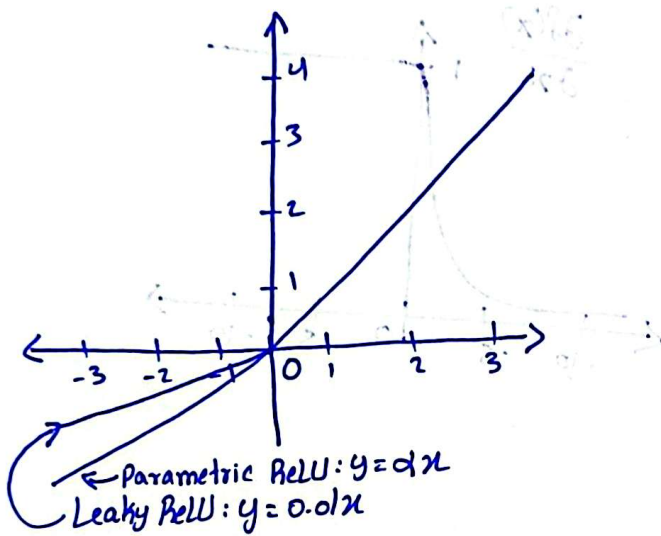
1

-6  -4  -2  2  4  6

$-\infty$ to $+\infty$

> Advantages

① All the advantages of ReLU

② Overcome the dead neuron issue

③ Less computational expensive compare to sigmoid and tanh

> Disadvantages

① It doesnot provide much learning for -ve +₂input as alpha value is constant

# ⟶ Parametric ReLU

$$\begin{cases} \alpha x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$$



← Parametric ReLU : $y = \alpha x$
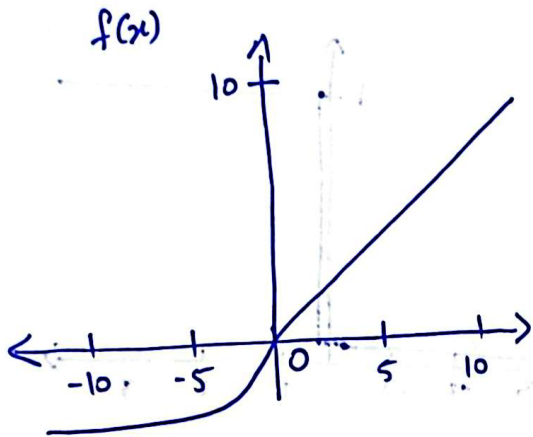Leaky ReLU : $y = 0.01 x$

$-\infty$ to $+\infty$

## > Advantages

① Slight advantage over leaky function as we can change learnable parameter alpha.

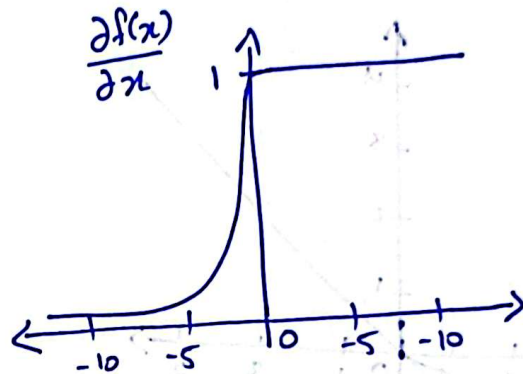② Less expensive compare to exponential activation function.

## > Disadvantages

① Does not provide more learning for -ve dataset

$\longrightarrow$ Exponential Linear Unit (ELU)

$$\begin{cases} \alpha(e^x - 1) & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$$
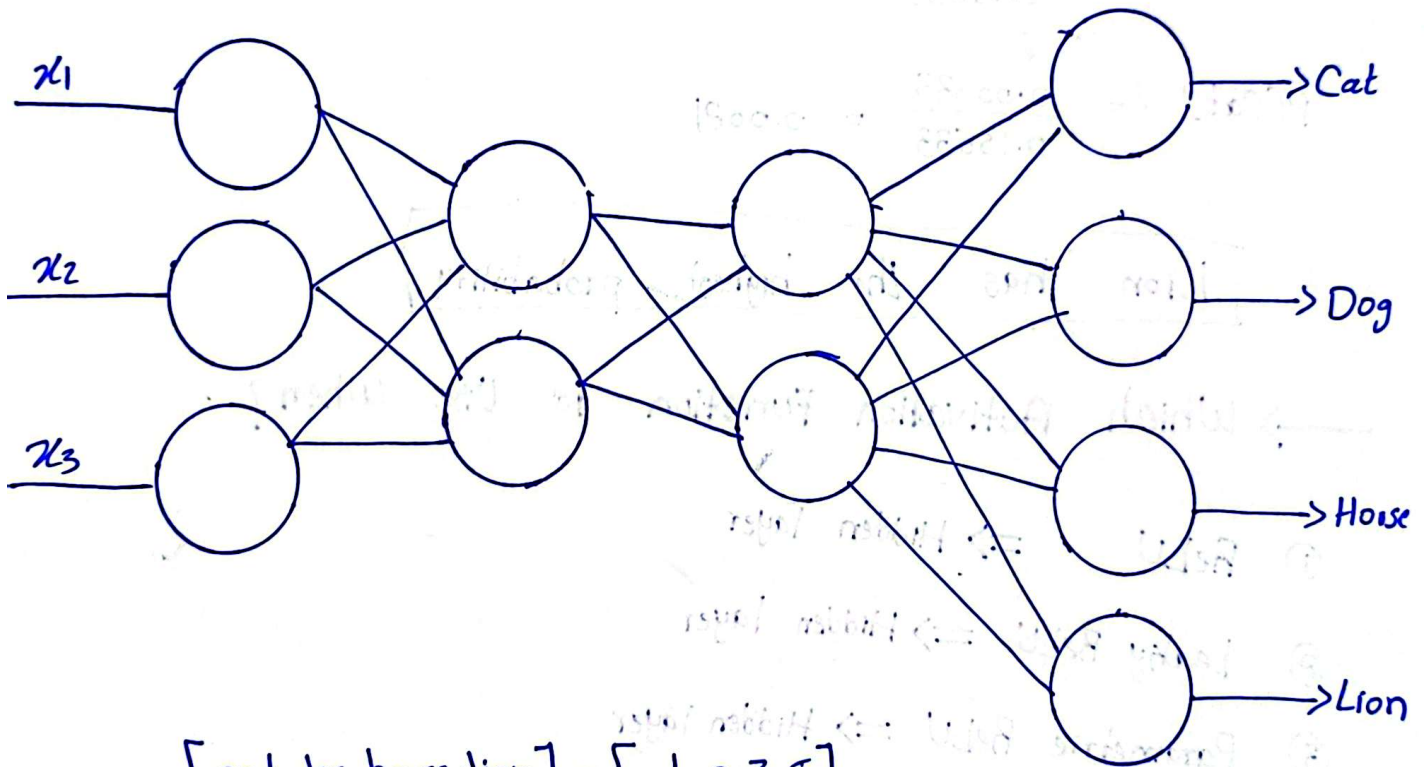


$-\alpha$ to $\infty$

> Advantages

① It is zero-centric

② Address issue of dying neuron

> Disadvantages

① Computationally expensive as exponential calculation is involved

# ⟶ Softmax Activation Function

Softmax activation function transforms the raw outputs of the neural network into a vector of probabilities, essentially a probability distribution over the input classes. Consider a multiclass classification problem with N classes.



$$[cat, dog, horse, lion] = [-1, 0, 3, 5]$$

$$Softmax = \frac{e^{y_i}}{\sum_{k=0}^{n} e^{y_k}} \quad ; \quad y_i = output \times weight + bias.$$

$$Cat = \frac{e^{-1}}{e^{-1+0+3+5}} = 0.003 \quad , \quad Dog = \frac{e^{0}}{e^{7}} = 0.0024$$

$$Horse = \frac{e^{3}}{e^{7}} = 0.0183 \quad , \quad Lion = \frac{e^{5}}{e^{7}} = 0.1353$$

$$P(Lion) = \frac{0.1353}{0.00033 + 0.0024 + 0.0183 + 0.1353} = 0.8654$$

$$P(Horse) = \frac{0.0183}{0.15633} = 0.11$$

$$P(Dog) = \frac{0.0024}{0.15633} = 0.015$$

$$P(Cat) = \frac{0.00033}{0.15633} = 0.0021$$

| Lion has the highest probability |
| --- |

⟶ Which Activation Function to Use when?

① ReLU ⟹ Hidden layer

② Leaky ReLU ⟹ Hidden layer

③ Parametric ReLU ⟹ Hidden layer

④ Exponential Linear Unit ⟹ Hidden layer

⑤ Sigmoid ⟹ Binary classification, output layer

⑥ Tanh ⟹ Hidden layer, shallow networks

⑦ Softmax ⟹ Multiclass classification, output layer